# PLFS Update

LANL: John Bent, HB Chen, David Gunter, Gary Grider, Sam Gutierrez, Adam Manzanares, Ben McClelland, Dave Montoya, James Nunez, Alfred Torrez, Meghan Wingate
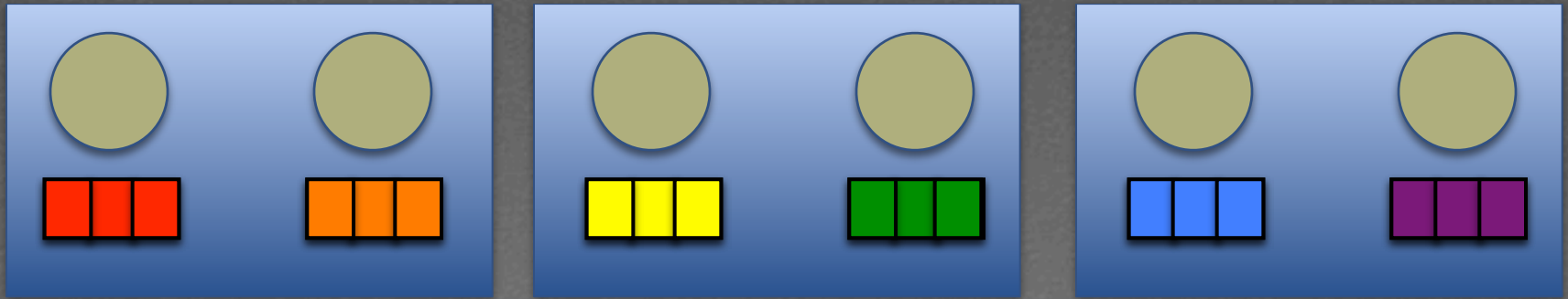
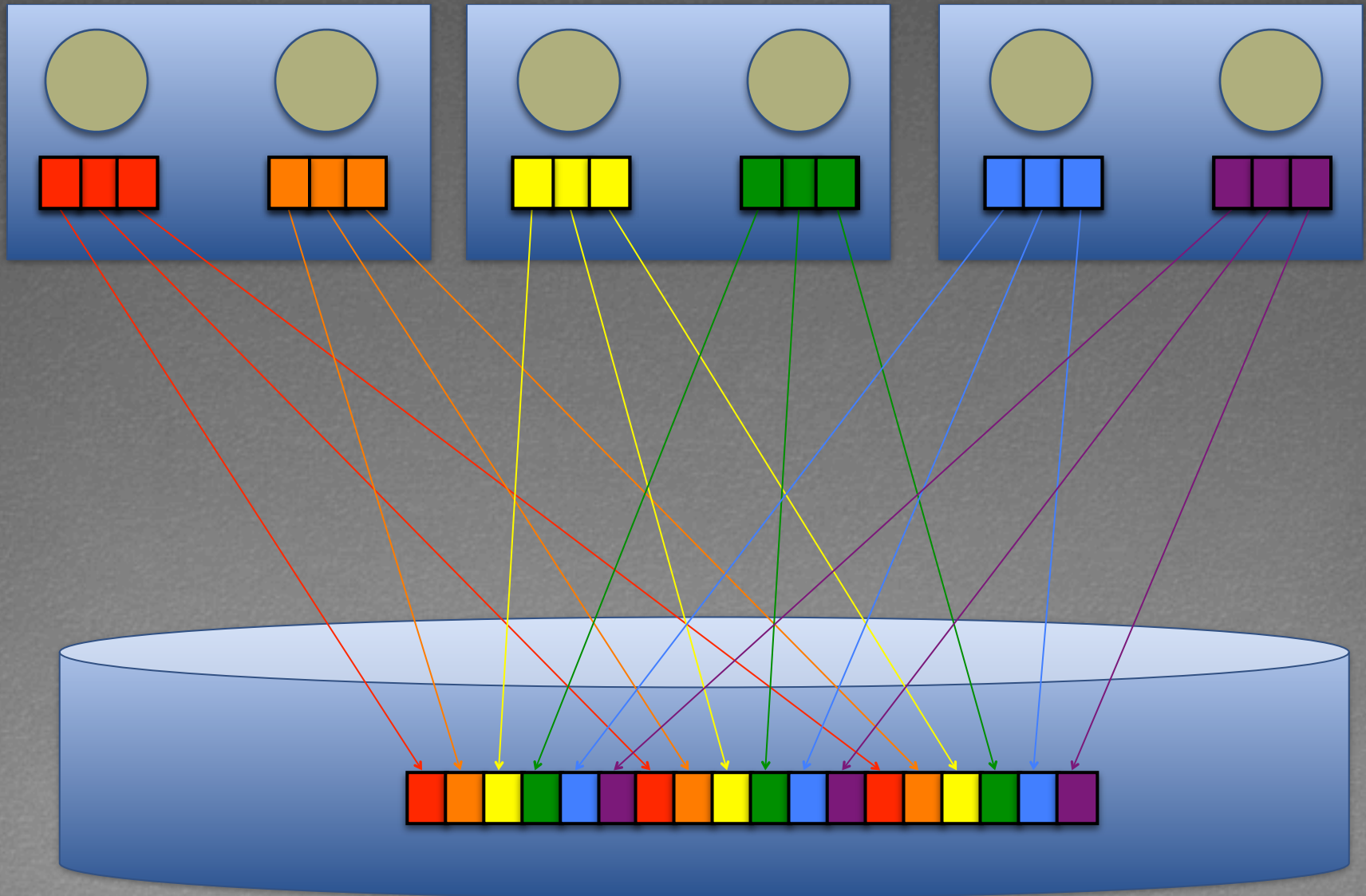CMU: Garth Gibson, Milo Polte

PSC: Paul Nowoczinski

# PLFS Overview

ଅ Virtual parallel file system

ଅ Interposes between application and another physical parallel file system
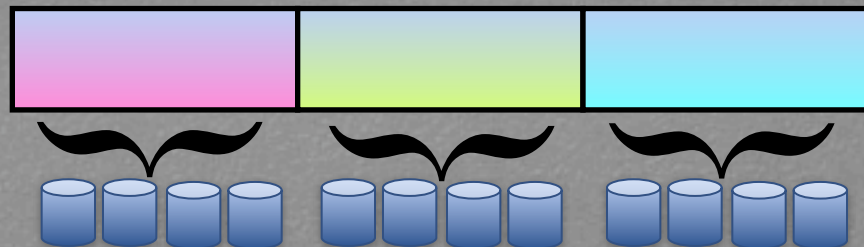
ଅ Reorganizes application IO
    ଅ From N-1 into N-N
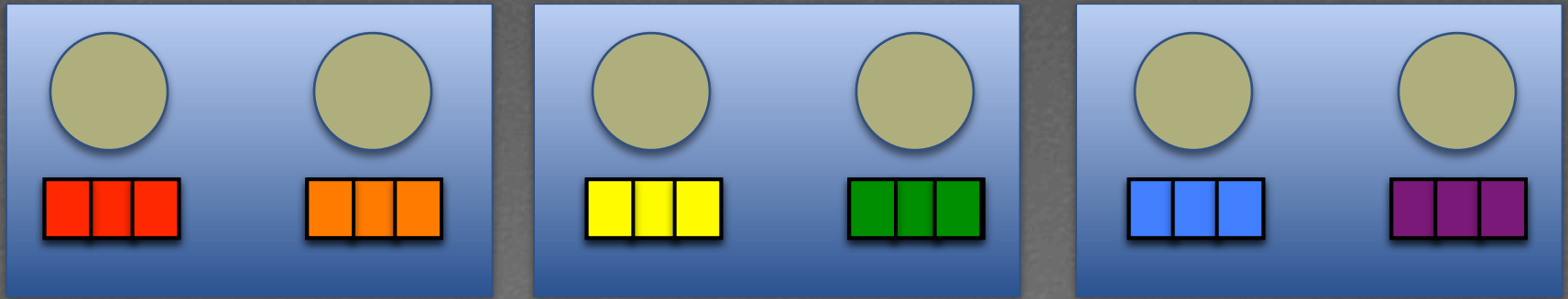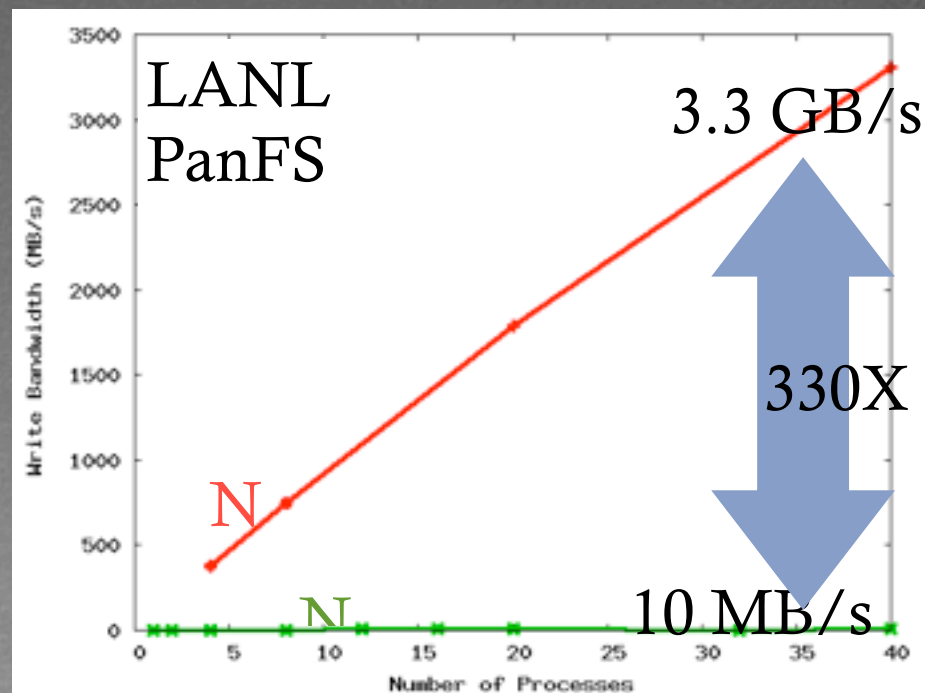
# Challenge of N-1 Strided

# Challenge of N-1 Strided

# Challenge of N-1 Strided

# N-1 Performance Penalty



LANL
GPFS

4.5 GB/s

45X

100 MB/s

PSC
Lustre

1.8 GB/s

90X

20 MB/s

LANL
PanFS

3.3 GB/s

330X

N

N

10 MB/s

N-N
N-1

# PLFS High Level



Physical Underlying Parallel File System

# PLFS High Level



PLFS Virtual Layer

Physical Underlying Parallel File System

# PLFS High Level



Physical Underlying Parallel File System

# PLFS High Level



Physical Underlying Parallel File System

PLFS Checkpoint BW Summary

# PLFS Status:
# Production Version 1.0 imminent

ଔ Current large push at LANL into production

    ଔ Running on Roadrunner for friendly users

    ଔ Big debug push to get onto Roadrunner for all users

    ଔ And other supercomputers as well: current and future

ଔ New features needed for Version 1.0 production release

    ଔ MPI-IO interface

    ଔ Thread-pools for archiving workloads

# PLFS Stack

# Write-optimized, now discovering read issues

&#8478;  Read open time for files created by lots of writers

&#8478;  Large sequential read bandwidth for files created with lots of small writes

# open() for read



Application
 fd = open(foo, O_RDONLY)
· · · · · · · · · · · · · · · · · · · · · · · · · ·
PLFS
    index = new Index()
    foreach index chunk c:
        index->add( c )

| Logical offset | Length | Chunk ID | Chunk Offset |
|---|---|---|---|
|  |  |  |  |
|  |  |  |  |
|  |  |  |  |
|  |  |  |  |
|  |  |  |  |
|  |  |  |  |

/foo/

/host1/

/host2/

/data.131  /data.132  /index

/data.279  /data.281  /index

**"PLFS Container"**

# open() for read

| Logical offset | Length | Chunk ID | Chunk Offset |
|---|---|---|---|
| | | | |
| | | | |
| | | | |
| | | | |
| | | | |
| | | | |

Application
 fd = open(foo, O_RDONLY)
· · · · · · · · · · · · · · · · · · · · · · · · · · · · ·
PLFS
    index = new Index()
    foreach index chunk c:
        index->add( c )

Replace foreach with threads

/foo/

/host1/                                    /host2/

/data.131   /data.132   /index      /data.279   /data.281   /index

**"PLFS Container"**

# Improving Read Open Times



**Read Open Times**

One Thread

Sixteen Threads

- plfs_threads=1
- plfs_threads=2
- plfs_threads=4
- plfs_threads=8
- plfs_threads=16

Read Open Time (s)

Initial Writers

Figure 1: experiment, right(target,3):read_file_open_wait_time_max, where mpihost like 'rrz' && user like 'johnbent' && description rlike 'hpss_plfs.reads5' && isnull(error) && read_file_open_wait_time_max < 20 && plfs_threads<=16 && total_size_mb > 20*1024 (1405 rows, 1280857844)

# Improving Read Bandwidths

**Application**
 read(fd, offset, len, buf)
· · · · · · · · · · · · · · · · · · · · · · · ·

**PLFS**
 foreach data chunk c:
  read(fd, offset+c,...)

↝ If file was created with small writes and read with large reads

↝ Each read may span multiple physical data chunks across multiple drives

/foo/

/host1/

/host2/

/data.131  /data.132   /index

/data.279  /data.281   /index

**"PLFS Container"**

# Improving Read Bandwidths

**Application**
read(fd, offset, len, buf)
. . . . . . . . . . . . . . . . . . . . . . . . . .

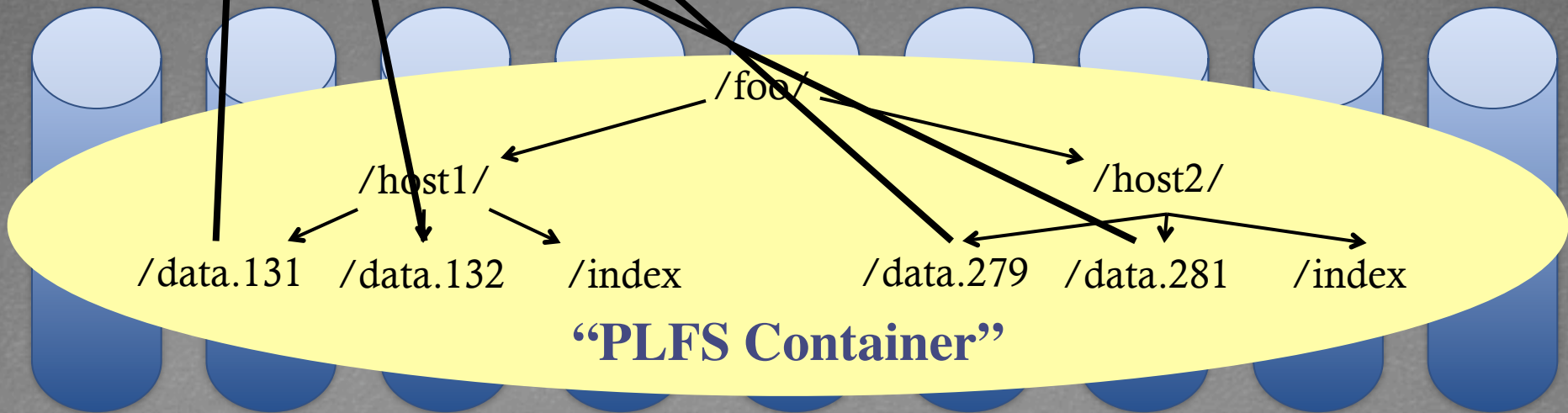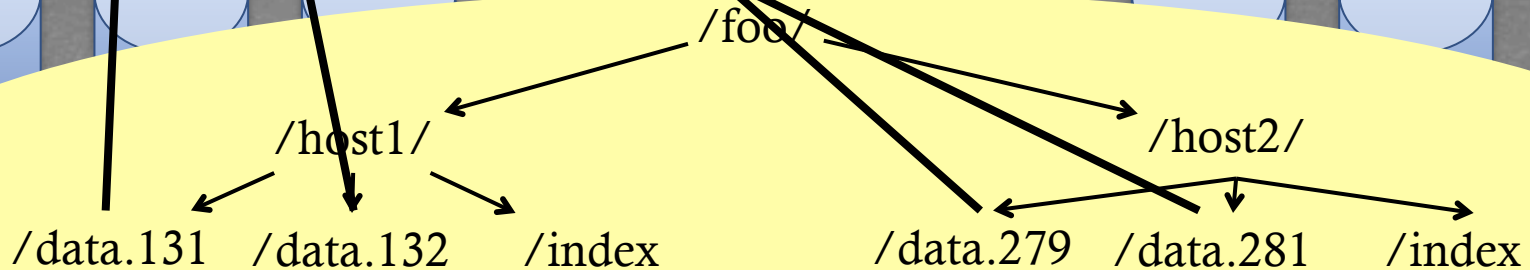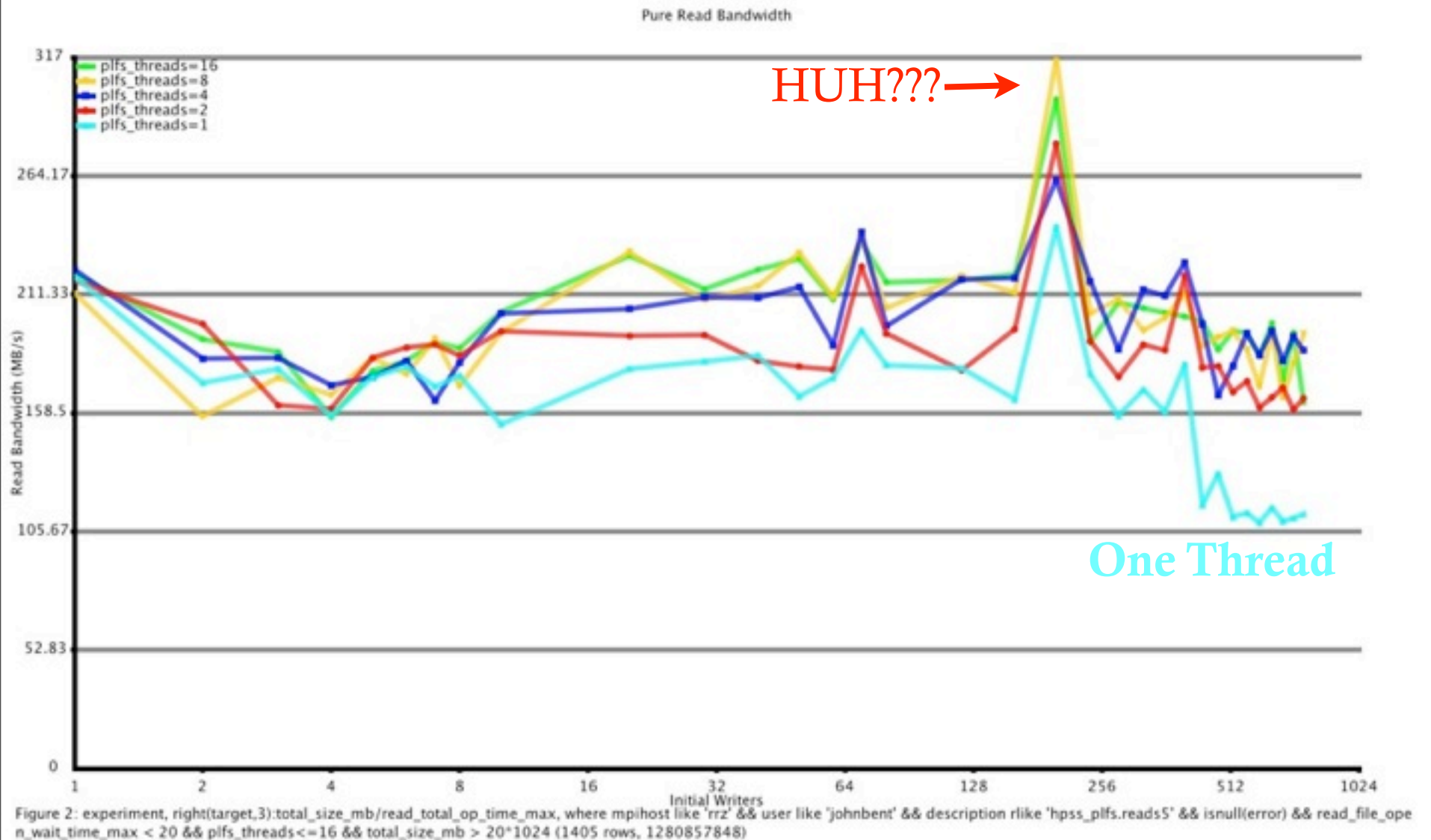**PLFS**
foreach data chunk c:
read(fd, offset+c,...)

ରେ If file was created with small writes and read with large reads

ରେ Each read may span multiple physical data chunks across multiple drives

Replace foreach with threads

/foo/

/host1/

/host2/

/data.131  /data.132  /index  /data.279  /data.281  /index

**"PLFS Container"**

# Multi-Threaded Reads



Figure 2: experiment, right(target,3):total_size_mb/read_total_op_time_max, where mpihost like 'rrz' && user like 'johnbent' && description rlike 'hpss_plfs.reads5' && isnull(error) && read_file_ope n_wait_time_max < 20 && plfs_threads<=16 && total_size_mb > 20*1024 (1405 rows, 1280857848)

Tuesday, August 3, 2010

# Conclusion

౭ Version 1.0 imminent

౭ Open source
  ౭ http://sourceforge.net/projects/plfs/
  ౭ > svn co https://plfs.svn.sourceforge.net/svnroot/plfs plfs

౭ Collaborations, contributions, bug reports very welcome


౭ Version 2.0 not imminent
  ౭ Further virtualization and re-organization of user data
  ౭ Distribute metadata workload transparently
  ౭ Use communicators in MPI to optimize PLFS ADIO layer